# TCSS 422: OPERATING SYSTEMS

**Three Easy Pieces:
Lock Based Data Structures,
Condition Variables**

Wes J. Lloyd
Institute of Technology
University of Washington - Tacoma

April 25, 2018      TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

---

## OBJECTIVES

- Assignment 1 – MASH Shell

- Lock Based Data Structures – Ch. 29

- Condition Variables – Ch. 30

- Quiz 3 – Lock-Based Data Structure Coding Activity

April 25, 2018      TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma      L8.2

---

## FEEDBACK – 4/23

- Do we have to know all the variations of Spin Locks?
  - Ch. 28 is just too much

- Basic spin lock
  - Polling / busy-waiting
  - When is it apprioriate?
- Test-and-set spin lock
- Compare-and-swap spin lock
  - *It is good to know what each successive version adds*

- Can you give us more practice questions on calculating the average response and turnaround time
  - Practice midterm next Monday 4/30

April 25, 2018      TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma      L8.3

---

## FEEDBACK – HW1

- We don't understand:
  "While writing your MASH program using processes, consider why it is non-trival to simply redirect the output stream of each fork command to a stream and use the C sscanf() function to consolidate/aggregate the output at the end... "

- *Comment is ambiguous*
- *In C (and Java) you could open arbitrary input and output streams that are not associated with a file on the disk. It would then be possible to redirect each exec command's output stream to streams not associated with files.*
- *Idea is to have "temporary in-memory buffers", in place of sending output to temporary files on the disk.*
- *I did not try this.*
- *For this assignment, it is easy enough to follow the example code and redirect exec output to temporary files:*

  *http://faculty.washington.edu/wlloyd/courses/tcss422/examples/exec2.c*

April 25, 2018      TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma      L8.4

---

## HW1 - 2

- We are executing commands fine and writing output. However, when it comes to the final process of displaying the output files in the correct order, it seems intuitive to use the Linux system() command.
- Is this acceptable?

- *This solution works, but it shouldn't be difficult to write a C routine that opens a file, reads it line-by-line, and displays output.*
- *This could be a generic, standalone routine.*

- *Example C should exist online by searching via Google to support accomplishing this.*

April 25, 2018      TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma      L8.5

---

## HW1 - 3

- For notification of process completion, the assignment states:
  First process finished...
  Second process finished...
  Third process finished...

- Should it always be that order? Or could it be:
  Third process finished...
  First process finished...
  Second process finished...

- *Reporting the order in which specific processes end is not required. Just report that processes *ARE* ending !*
- *There is no requirement to say which child process/PID finishes in what order, etc.*
- *FEATURE: Provides a notification message stating processes are finishing and work is proceeding.*

April 25, 2018      TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma      L8.6

---

---

# CHAPTER 29 – LOCK BASED DATA STRUCTURES

April 25, 2018
TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

L8.7

---

## LOCK-BASED CONCURRENT DATA STRUCTURES

- Adding locks to data structures make them **thread safe**.

- Considerations:
  - Correctness
  - Performance
  - Lock granularity

April 25, 2018
TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

L8.8

---

## CONCURRENT COUNTERS - PERFORMANCE

- iMac: four core Intel 2.7 GHz i5 CPU
- Each thread increments counter 1,000,000 times

Traditional vs. sloppy counter
Sloppy Threshold (S) = 1024

**Synchronized counter scales poorly.**

April 25, 2018
TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

L8.9

---

## PERFECT SCALING

- Achieve (N) performance gain with (N) additional resources

- Throughput:
- Transactions per second

- 1 thread/core
- N = 100 tps

- 10 threads/cores
- N = 1000 tps

April 25, 2018
TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

L8.10

---

## SLOPPY COUNTER

- Provides single logical shared counter
  - Implemented using local counters for each ~CPU core
    - 4 CPU cores = 4 local counters & 1 global counter
    - Local counters are synchronized via local locks
  - Global counter is updated periodically
    - Global counter has lock to protect global counter value
    - Sloppiness threshold (S):
      Update threshold of global counter with local values
    - Small (S): more updates, more overhead
    - Large (S): fewer updates, more performant, less synchronized
- Why this implementation?
  Why do we want counters local to each CPU Core?

April 25, 2018
TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

L8.11

---

## SLOPPY COUNTER - 2

- Update threshold (S) = 5
- Synchronized across four CPU cores
- Threads update local CPU counters

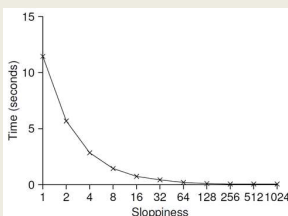| Time | $L_1$ | $L_2$ | $L_3$ | $L_4$ | G |
|------|------|------|------|------|---|
| 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 1 | 1 | 0 |
| 2 | 1 | 0 | 2 | 1 | 0 |
| 3 | 2 | 0 | 3 | 1 | 0 |
| 4 | 3 | 0 | 3 | 2 | 0 |
| 5 | 4 | 1 | 3 | 3 | 0 |
| 6 | 5 → 0 | 1 | 3 | 4 | 5 (from $L_1$) |
| 7 | 0 | 2 | 4 | 5 → 0 | 10 (from $L_4$) |

April 25, 2018
TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma

L8.12

---

## THRESHOLD VALUE S

- Consider 4 threads increment a counter 1000000 times each
- Low $S$ → What is the consequence?
- High $S$ → What is the consequence?

## SLOPPY COUNTER - EXAMPLE

- Example implementation

- Also with CPU affinity

## CONCURRENT LINKED LIST - 1

- Simplification - only basic list operations shown
- Structs and initialization:

```
1       // basic node structure
2       typedef struct __node_t {
3               int key;
4               struct __node_t *next;
5       } node_t;
6
7       // basic list structure (one used per list)
8       typedef struct __list_t {
9               node_t *head;
10              pthread_mutex_t lock;
11      } list_t;
12
13      void List_Init(list_t *L) {
14              L->head = NULL;
15              pthread_mutex_init(&L->lock, NULL);
16      }
17
(Cont.)
```

## CONCURRENT LINKED LIST - 2

- Insert – adds item to list
- Everything is critical!
  - There are two unlocks

```
(Cont.)
18      int List_Insert(list_t *L, int key) {
19              pthread_mutex_lock(&L->lock);
20              node_t *new = malloc(sizeof(node_t));
21              if (new == NULL) {
22                      perror("malloc");
23                      pthread_mutex_unlock(&L->lock);
24              return -1; // fail
26              new->key = key;
27              new->next = L->head;
28              L->head = new;
29              pthread_mutex_unlock(&L->lock);
30              return 0; // success
31      }
(Cont.)
```

## CONCURRENT LINKED LIST - 3

- Lookup – checks list for existence of item with key
- Once again everything is critical
  - Note - there are also two unlocks

```
(Cont.)
32
32      int List_Lookup(list_t *L, int key) {
33              pthread_mutex_lock(&L->lock);
34              node_t *curr = L->head;
35              while (curr) {
36                      if (curr->key == key) {
37                              pthread_mutex_unlock(&L->lock);
38                              return 0; // success
39                      }
40                      curr = curr->next;
41              }
42              pthread_mutex_unlock(&L->lock);
43              return -1; // failure
44      }
```

## CONCURRENT LINKED LIST

- First Implementation:
  - Lock **everything** inside Insert() and Lookup()
  - If malloc() fails lock must be released
    - Research has shown "***exception-based control flow***" to be error prone
    - 40% of Linux OS bugs occur in rarely taken code paths
    - Unlocking in an exception handler is considered a poor coding practice
    - There is nothing specifically wrong with this example however

- Second Implementation …

---

## CCL – SECOND IMPLEMENTATION

- Init and Insert

```
1     void List_Init(list_t *L) {
2         L->head = NULL;
3         pthread_mutex_init(&L->lock, NULL);
4     }
5
6     void List_Insert(list_t *L, int key) {
7         // synchronization not needed
8         node_t *new = malloc(sizeof(node_t));
9         if (new == NULL) {
10            perror("malloc");
11            return;
12        }
13        new->key = key;
14
15        // just lock critical section
16        pthread_mutex_lock(&L->lock);
17        new->next = L->head;
18        L->head = new;
19        pthread_mutex_unlock(&L->lock);
20    }
21
```

| April 25, 2018 | TCSS422: Operating Systems [Spring 2018] Institute of Technology, University of Washington - Tacoma | L8.19 |

---

## CCL – SECOND IMPLEMENTATION - 2

- Lookup

```
(Cont.)
22    int List_Lookup(list_t *L, int key) {
23        int rv = -1;
24        pthread_mutex_lock(&L->lock);
25        node_t *curr = L->head;
26        while (curr) {
27            if (curr->key == key) {
28                rv = 0;
29                break;
30            }
31            curr = curr->next;
32        }
33        pthread_mutex_unlock(&L->lock);
34        return rv; // now both success and failure
35    }
```

| April 25, 2018 | TCSS422: Operating Systems [Spring 2018] Institute of Technology, University of Washington - Tacoma | L8.20 |

---

## CONCURRENT LINKED LIST PERFORMANCE

- Using a single lock for entire list is not very performant
- Users must "wait" in line for a single lock to access/modify any item
- Hand-over-hand-locking (lock coupling)
  - Introduce a lock for each node of a list
  - Traversal involves handing over previous node's lock, acquiring the next node's lock…
  - Improves lock granularity
  - Degrades traversal performance

- Consider hybrid approach
  - Fewer locks, but more than 1
  - Best lock-to-node distribution?

| April 25, 2018 | TCSS422: Operating Systems [Spring 2018] Institute of Technology, University of Washington - Tacoma | L8.21 |

---

## MICHAEL AND SCOTT CONCURRENT QUEUES

- Improvement beyond a single master lock for a queue (FIFO)
- Two locks:
  - One for the **head** of the queue
  - One for the **tail**
- Synchronize enqueue and dequeue operations

- Add a dummy node
  - Allocated in the queue initialization routine
  - Supports separation of head and tail operations

- Items can be added and removed by separate threads at the same time

| April 25, 2018 | TCSS422: Operating Systems [Spring 2018] Institute of Technology, University of Washington - Tacoma | L8.22 |

---

## CONCURRENT QUEUE

- Remove from queue

```
1     typedef struct __node_t {
2         int value;
3         struct __node_t *next;
4     } node_t;
5
6     typedef struct __queue_t {
7         node_t *head;
8         node_t *tail;
9         pthread_mutex_t headLock;
10        pthread_mutex_t tailLock;
11    } queue_t;
12
13    void Queue_Init(queue_t *q) {
14        node_t *tmp = malloc(sizeof(node_t));
15        tmp->next = NULL;
16        q->head = q->tail = tmp;
17        pthread_mutex_init(&q->headLock, NULL);
18        pthread_mutex_init(&q->tailLock, NULL);
19    }
20
(Cont.)
```

| April 25, 2018 | TCSS422: Operating Systems [Spring 2018] Institute of Technology, University of Washington - Tacoma | L8.23 |

---

## CONCURRENT QUEUE - 2

- Add to queue

```
(Cont.)
21    void Queue_Enqueue(queue_t *q, int value) {
22        node_t *tmp = malloc(sizeof(node_t));
23        assert(tmp != NULL);
24
25        tmp->value = value;
26        tmp->next = NULL;
27
28        pthread_mutex_lock(&q->tailLock);
29        q->tail->next = tmp;
30        q->tail = tmp;
31        pthread_mutex_unlock(&q->tailLock);
32    }
(Cont.)
```

| April 25, 2018 | TCSS422: Operating Systems [Spring 2018] Institute of Technology, University of Washington - Tacoma | L8.24 |

## CONCURRENT HASH TABLE

- Consider a simple hash table
  - Fixed (static) size
  - Hash maps to a bucket
    - Bucket is implemented using a concurrent linked list
    - One lock per hash (bucket)
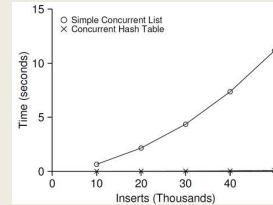    - Hash bucket is a linked lists

April 25, 2018 — TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma — L8.25

## INSERT PERFORMANCE – CONCURRENT HASH TABLE

- Four threads – 10,000 to 50,000 inserts
  - iMac with four-core Intel 2.7 GHz CPU

The simple concurrent hash table scales magnificently

April 25, 2018 — TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma — L8.26

## CONCURRENT HASH TABLE

```
1    #define BUCKETS (101)
2
3    typedef struct __hash_t {
4            list_t lists[BUCKETS];
5    } hash_t;
6
7    void Hash_Init(hash_t *H) {
8            int i;
9            for (i = 0; i < BUCKETS; i++) {
10                   List_Init(&H->lists[i]);
11           }
12   }
13
14   int Hash_Insert(hash_t *H, int key) {
15           int bucket = key % BUCKETS;
16           return List_Insert(&H->lists[bucket], key);
17   }
18
19   int Hash_Lookup(hash_t *H, int key) {
20           int bucket = key % BUCKETS;
21           return List_Lookup(&H->lists[bucket], key);
22   }
```

April 25, 2018 — TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma — L8.27

## LOCK-FREE DATA STRUCTURES

- Lock-free data structures in Java

- Java.util.concurrent.atomic package
- Classes:
  - AtomicBoolean
  - AtomicInteger
  - AtomicIntegerArray
  - AtomicIntegerFieldUpdater
  - AtomicLong
  - AtomicLongArray
  - AtomicLongFieldUpdater
  - AtomicReference

- See: https://docs.oracle.com/javase/7/docs/api/java/util/concurrent/atomic/package-summary.html

April 25, 2018 — TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma — L8.28

# CHAPTER 30 – CONDITION VARIABLES

April 25, 2018 — TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma — L8.29

## CONDITION VARIABLES

- There are many cases where a thread wants to wait for another thread before proceeding with execution

- Consider when a precondition must be fulfilled before it is meaningful to proceed …

April 25, 2018 — TCSS422: Operating Systems [Spring 2018]
Institute of Technology, University of Washington - Tacoma — L8.30

# QUESTIONS